

What We Found on Our Way to Building a Classifier: A Critical Analysis of the AHA Screening Questionnaire

Quazi Abidur Rahman¹, Sivajothi Kanagalingam², Aurelio Pinheiro²,
Theodore Abraham², and Hagit Shatkay^{1,3}

¹Computational Biology and Machine Learning Lab, School of Computing, Queen's University,
Kingston, ON, Canada
quazi@cs.queensu.ca

²Heart and Vascular Institute, Johns Hopkins University, Baltimore, MD, USA
kanagalingam.jothi@gmail.com, {apinhei5, tabraha3}@jhmi.edu

³Dept. of Computer and Information Sciences & Center for Bioinformatics and Computational
Biology, University of Delaware, Newark, DE, USA
shatkay@cis.udel.edu

Abstract. The American Heart Association (AHA) has recommended a 12-element questionnaire for pre-participation screening of athletes, in order to reduce and hopefully prevent sudden cardiac death in young athletes. This screening procedure is widely used throughout the United States, but its efficacy for discriminating *Normal* from *Non-normal* heart condition is unclear. As part of a larger study on cardiovascular disorders in young athletes, we set out to train machine-learning-based classifiers to automatically categorize athletes into risk-levels based on their answers to the AHA-questionnaire. We also conducted information-based and probabilistic analysis of each question to identify the ones that may best predict athletes' heart condition. *However, surprisingly*, the results indicate that the AHA-recommended screening procedure itself does not effectively distinguish between *Normal* and *Non-normal* heart as identified by cardiologists using Electro- and Echo-cardiogram examinations. Our results suggest that ECG and Echo, rather than the questionnaire, should be considered for screening young athletes.

1 Introduction

Inherited cardiovascular disease is the main cause of sudden cardiac death (SCD) in young athletes. In the United States the incidence has been reported as 1:50,000 – 1:100,000 per year [1–3]. A larger study in the Veneto region in Italy reported an incidence rate of SCD of 2.1 per 100,000 athletes annually as a result of cardiovascular disease [1]. While the incidence of SCD is lower in comparison to other causes of death, it is disconcerting in that these deaths occur in young and otherwise perceived-to-be healthy individuals, most often without any prior cardiac symptoms. Moreover, as most of these deaths occur in athletes of high-school age [1,4], they are a cause for much concern in the media, the public and the medical community.

Initial screening through electrocardiogram (ECG) and echocardiogram (Echo) is a first step for identifying morphological anomalies that can lead to cardiac abnormaliti-

Table 1. The AHA 12-element Screening Guidelines [8]

Guideline #	Question Type	Question Contents as described in the AHA guideline
1	Personal History	Exertional chest pain/discomfort?
2		Unexplained syncope/near-syncope?
3		Excessive exertional and unexplained dyspnea/fatigue, associated with exercise?
4		Prior recognition of a heart murmur?
5		Elevated systemic blood pressure ?
6	Family History	Premature death (sudden and unexpected, or otherwise) before age 50 years due to heart disease, in at least one relative?
7		Disability from heart disease in a close relative younger than 50 years of age?
8		Specific knowledge of certain cardiac conditions in family members: hypertrophic or dilated cardiomyopathy, long-QT syndrome or other ion channelopathies, Marfan syndrome, or clinically important arrhythmias?
9	Physical Exam	Heart murmur
10		Femoral pulses to exclude aortic coarctation
11		Physical stigmata of Marfan syndrome
12		Brachial artery blood pressure (sitting position)

es, and in extreme cases to sudden death. However, due to considerations involving speed, ease of administration and cost, these standard procedures, while often used in Europe [5] are not used for large-scale screening of young athletes in the United States. As an alternative preventive measure, the American Heart Association (AHA) has recommended a screening procedure [6], intended as a cost-effective, practical initial measure for pre-participation screening of athletes. In the United States, the use of this screening procedure has steadily increased over the years since 1997 [7].

The current, revised, AHA pre-participation screening recommendations were published in 2007, and include 12-element screening guidelines [8] (see Table 1). Under these guidelines, each athlete answers several questions concerning personal and family history and undergoes a physical examination (we refer to the combination of questions and physical exam as “*the questionnaire*”). If any of the questions is answered in the affirmative, or if the physical examination suggests an abnormality, the athlete is then referred for a more extensive cardiologic evaluation through ECG and Echo, in which responses that are *Non-normal* (i.e., deviate from the Normal measures established for athletes, but not conclusively abnormal) can be identified; athletes with *Non-normal* results are referred for further, more extensive testing to verify whether any serious heart condition is present. A preliminary study by our group [9] (presented as an abstract at the AHA symposium), has broadly suggested low predictive power of the AHA screening procedure, without considering its explicit elements and their predictive value.

As a component within a large-scale research of adverse heart conditions, which extensively studies the efficacy of the questionnaire and its possible contribution to predicting cardiac irregularities, we set out to pursue what appeared to be a straightforward task: namely, training a machine-learning-based classifier, based on the answers to the questionnaire from several hundred athletes, in order to automatically

predict from these answers the athletes' heart condition. The "heart condition" for the purpose of this study was either *Normal* or *Non-normal*, as determined by a cardiologist based on ECG and Echo readings. The cardiologist's adjudication, which is based solely on ECG and Echo, serves here as the "gold-standard" to which the AHA guidelines results are compared. We expected to be able to effectively train such a classifier from the questionnaire data, due to the hypothesis driving the AHA guidelines as discussed above: namely, that the answers to the pre-screening questionnaire can indeed be correlated with the diagnosis obtained from the more extensive and time-consuming, Echo and ECG tests, administered by a physician. Intending to follow the common machine-learning procedures for learning a classifier from data (e.g., [10]) (we also aimed to select the most informative features, that is, identify the items in the AHA-based pre-screening procedure, whose answers are the most predictive of the cardiologist's adjudication. Machine learning methods have been widely used for disease prediction, risk assessment and patient classification. For instance, in the field of cardiology, arrhythmia classification was performed using support vector machines [11, 12], linear discriminant analysis [13] and artificial neural networks [14]. As another example, naïve Bayes classifiers have been used for diagnosis and risk assessment of Long-QT syndrome in children from clinical data [15]. In the area of cancer diagnosis and prediction, methods such as support vector machines [16], logistic regression [17] and random forests [18] have been applied. We thus anticipated that by using filled-in questionnaires from a relatively large population of young athletes, we could train a classifier to distinguish between athletes with potential cardiovascular abnormalities (as determined by ECG and Echo tests) from normal ones.

Notably, the screening through the AHA questionnaire is intended as a means to avoid the more costly and cumbersome Echo and ECG tests. Thus the underlying assumption in administering the AHA procedure is that athletes who require further screening (those whose ECG or Echo would thus not be completely *Normal*) would indeed be identified in the screening and referred for further examination (ECG, Echo - and if needed even more extensive testing), while athletes who do not need further screening would have their questions and basic physical show completely normal answers. Based on this insight, the expectation was that the answers to the questionnaire should be predictive of the Echo/ECG results. As such, our original goal was to train a machine-learning-based classifier that will take as input the results obtained from the screening based on the 12-element AHA guidelines for each athlete and predict the cardiologist's Echo/ECG-based adjudication. In this study we rigorously apply classification techniques and investigate the information-content of each item in the questionnaire. We also conduct probabilistic analysis of the positive and negative answers and their correlation with ECG/Echo test results. However, the classification results and the information contents of the different items, as well as the results from the probabilistic analysis, expose significant shortcomings in the pre-screening procedure itself. Thus, what started as a classification task, ended up as an in-depth informatics-driven analysis, revealing important issues with the AHA screening procedure, whose use is advocated as the primary screening tool for athletes.

While this article begins by discussing what appears to be a negative result, its main contribution and the significance of the presented research lies in employing the sa-

me statistical, information-based methods that are typically used for developing diagnostic/predictive machine-learning tools, to effectively expose important shortcomings in the current screening procedure. It also points out that other, more discerning, procedures may be required for effective pre-participation screening of athletes (at least until a questionnaire is devised with better predictive capability). Hence, our results suggest that ECG and (possibly Echo) should be considered for screening athletes in the United States. We note that ECG is being used for screening of athletes in Europe, especially in Italy [5] and has been recommended by the consensus statement of the European Society of Cardiology [1].

Throughout the rest of the paper we describe the AHA-based questionnaire data, the analysis applied, and the operative conclusions, suggesting that the questionnaire is not an effective tool for assessing risk in young athletes, and that alternative procedures need to be considered.

2 Data Used in this Study

The study included 470 participants, all of whom are young athletes participating at state-level athletic events. They were all asked to fill a questionnaire consisting of 12 *Yes/No* questions as shown in Table 2 (*Q1-Q12*), corresponding to AHA elements 1-8 shown in Table 1. They have also undergone a standard, basic physical exam corresponding to AHA elements 9-12 in Table 1. The results of the physical (which can either be normal or abnormal), are listed as Question 13 (*Q13*) in Table 2. Notably, the AHA 12-elements are intended to be clear to physicians but not necessarily to laymen. Therefore, the questionnaire filled by the athletes, as shown in Table 2, uses simply-phrases questions that correspond to each element's intention. In several cases more than one question is needed to cover an element, and some questions address more than a single element. The element number(s) covered by each question is shown in the rightmost column of Table 2.

In addition to answering questions *Q1-Q12* and undergoing the basic physical (*Q13*), the participants have separately undergone ECG and Echo tests. The latter two tests were evaluated by an expert cardiologist to draw a more conclusive adjudication regarding each individual's heart condition, based on measurable, observable cardiac parameters as opposed to questions. The two possible conclusions were: *Normal* and *Non-normal*, where *Non-normal* heart condition means that further extensive cardiological evaluation of the athlete is required. The cardiologist's adjudication was based solely on the ECG and Echo tests, and did not include any analysis or consideration of the questionnaire results. Of the 470 participants, 348 were categorized by the cardiologist as *Normal*, while 122 were categorized as *Non-normal*.

As not all participants answered all the questions, when analyzing individual questions for information content and conditional probabilities (Sections 3.2 and 3.3 below), we consider, per-question, only the number of answers that the question has actually received. In Section 3.1, we describe how the missing values are handled by the classifiers. The second row in Table 3 shows how many answers were received for each of the questions, while the third and fourth rows indicate how many of the answers were positive and how many of them were negative, respectively.

Table 2. The list of questions used in the questionnaire presented to the athletes in this study, along with the AHA guideline number to which each question corresponds

Quest. #	Question content as presented to athlete	AHA Guideline #
Q1	Dizziness/Passed Out during/after exercise?	2
Q2	Chest Pains or shortness of breath?	1
Q3	Become tired quicker than peers during exercise?	3
Q4	Heart murmur/disease?	4
Q5	Skipped heartbeats or racing heartbeats?	1 (discomfort), 4
Q6	Heart disease development or related death in family?	6
Q7	Does anyone in the family have fainting episodes or seizures?	6,7
Q8	Chest discomfort when active?	1
Q9	Have you been told you have high blood pressure?	5
Q10	Have you experiences seizures or exercise related asthma?	1,2
Q11	Anyone in family experienced heart surgery or have a pace-maker or defibrillator under the age of 50 years?	7
Q12	Anyone in family diagnosed with Cardiomyopathy, aneurysm, Marfan's, IHSS?	8
Q13	Physical examination results <i>abnormal</i> ?	9-12

3 Methods and Tools

Our analysis of the AHA questionnaire data started by applying classifiers to the data, and was followed by an information-content analysis of each question. We also performed probabilistic analysis of the answers to each question. These methods and related tools are presented in the following subsections.

3.1 The Classifiers

As a baseline for examining the feasibility of predicting the heart condition of young athletes using the AHA questions and physical examination as attributes, we applied three standard classification methods: naïve Bayes (e.g., [10]), random forests [19] and support vector machine (SVM) [20]. We used the standard classification packages in WEKA [21] for all three classifiers. The Random forests algorithm was implemented with 100 trees. SVM used Gaussian radial basis function as kernel¹, where the soft margin parameter C and the kernel parameter γ were selected after trying several combinations of the parameters and choosing the best one in terms of overall accuracy. To train/test and evaluate the performance of the classifiers, we used the standard 10-fold cross-validation procedure.

As not all participants answered all the questions, some values are missing in the questionnaires, as shown in Table 3. For classification purposes, we denote each miss-

¹ We have also tried linear kernel, but Gaussian radial basis kernel performed marginally better than the linear kernel.

Table 3. Number of answers received for each question, along with the number of positive and negative answers

	<i>Q1</i>	<i>Q2</i>	<i>Q3</i>	<i>Q4</i>	<i>Q5</i>	<i>Q6</i>	<i>Q7</i>	<i>Q8</i>	<i>Q9</i>	<i>Q10</i>	<i>Q11</i>	<i>Q12</i>	<i>Q13</i>
# of answers	469	466	466	436	431	380	423	466	440	468	459	367	451
# of positive answers	94	121	51	33	22	40	45	55	26	65	6	12	40
# of negative answers	375	345	415	403	409	340	378	411	414	403	453	355	411

ing value as *Not Known (NK)*. Hence, each athlete’s response to the questionnaire is represented as a 13-dimensional vector $(a_1, a_2, a_3, \dots, a_{13})$, where $a_i \in \{No, Yes, NK\}$, denoting a negative, a positive or a *Not Known* answer, respectively, to question Q_i . The intended task for each classifier is to assign each such instance (athlete) into one of the two possible classes: *Normal* or *Non-normal*. For the purpose of this study, the gold-standard, true class for each of the 470 athletes is as assigned by the cardiologist based on the results of the ECG and Echo tests (348 have *Normal* conclusion and 122 have *Non-normal* conclusion).

As the dataset is biased toward the *Normal* class, to correct for the imbalance, we used the procedure of sub-sampling from the over-represented class to create a balanced dataset for training/testing. Under the sub-sampling method, instances are chosen at random from the majority class to make the size of the two classes equal. By randomly selecting 122 instances from the *Normal* class and taking the whole subset of 122 *Non-normal* instances we obtain a balanced dataset. We have repeated the sub-sampling procedure 5 times to ensure stability of the results. The classifiers have been trained and tested on both the original and the balanced dataset. To evaluate the performance of the classifiers, we have used several standard measures, namely, the *Accuracy* (the proportion of correctly classified instances), as well as the widely used measures of *Recall* (Sensitivity), *Precision* (counterpart of Specificity), and *F-measure*. *Accuracy*, *Precision* and *Recall* are defined below, where true positives, denote *Non-normal* cases that are correctly classified as *Non-normal*:

$$Accuracy = (\# \text{ of correctly classified instances}) / (\text{Total number of instances}) ;$$

$$Precision = (\# \text{ of true positives}) / (\# \text{ of true positives} + \# \text{ of false positives}) ;$$

$$Recall = (\# \text{ of true positives}) / (\# \text{ of true positives} + \# \text{ of false negatives}).$$

The *F-measure* is the harmonic mean of the *Precision* and the *Recall*. The definition of the *F-measure* is: $F - \text{measure} = 2 \cdot (Precision \cdot Recall) / (Precision + Recall)$.

3.2 Information Content Analysis

As discussed in more detail in Section 4, using all the questions as attributes results in poor classification performance. Hence we investigated each question individually to assess its predictive capability. To measure each question’s predictive capability, we use the well-known Information Gain criterion (e.g., [10]). The information gain, calculated for each question, measures how much information is gained about the conclusion (*Normal* or *Non-normal*) when the answer to that question is obtained. It thus indicates how predictive the answer to a question is in classifying participants as having a *Normal* or a *Non-normal* heart-condition. It is calculated as the difference between the unconditional entropy associated with the conclusion and the conditional entropy of the conclusion given the answer to a question. These measures are formal-

ly defined as follows: Let C be the set of conclusions (class labels) and A_Q be the answer to question Q . The maximum likelihood estimate for the probability of the conclusion being *Normal*, or *Nor* for short, $Pr(C = Nor)$, is calculated as:

$$Pr(C = Nor) \approx (\# \text{ of participants with Normal concl.}) / (\text{Total \# of participants}),$$

while the probability of *Non-normal* (denoted *NNor*) conclusion is calculated as:

$$Pr(C = NNor) = 1 - Pr(C = Nor).$$

Similarly, we define the conditional probability of the conclusion to be *Normal* (or *Non-normal*) given the answer (*Yes* or *No*) to question Q . We define this probability, for a question Q , as: $Pr(C = W|A_Q = X)$ where W is either *Nor* or *NNor* and X is either *Yes* or *No*. The conditional probabilities are estimated from the observed proportions; e.g., the probability of the conclusion being *Non-normal* given that the answer for question Q is positive, $Pr(C = NNor|A_Q = Yes)$, is estimated as:

$$Pr(C = NNor|A_Q = Yes) \approx \frac{\# \text{ of participants with Non-normal conclusion and positive answer to } Q}{\text{Total \# of participants who have answered positively to } Q}.$$

The entropy of the conclusion, $H(C)$, is defined as:

$$H(C) = -[Pr(C = Nor) \log_2 Pr(C = Nor) + Pr(C = NNor) \log_2 Pr(C = NNor)].$$

Let the conditional entropy of the conclusion, given a positive or a negative answer be $H(C|A_Q = Yes)$ and $H(C|A_Q = No)$, respectively. The conditional entropy of the conclusions set C given the answer to a question Q is calculated as:

$$H(C|A_Q) = [Pr(A_Q = Yes) * H(C|A_Q = Yes) + Pr(A_Q = No) * H(C|A_Q = No)]$$

The information gain associated with question Q , $IG(C, A_Q)$, is formally defined as:

$$IG(C, A_Q) = H(C) - H(C|A_Q).$$

3.3 Probabilistic Analysis of the Questions

As all questions lead to a very low information gain (see Section 4), we investigated for each question whether a positive answer to it has a significantly higher probability of indicating *Non-normal* conclusion, compared to a negative answer. Any such question is expected to at least indicate a likely *Non-normal* conclusion (even if it does not reliably identify *Normal* conclusions). We note that correctly identifying *Non-normal* conclusion is more important than correctly predicting *Normal* conclusion, because failure to identify an athlete with a *Non-normal* conclusion can be potentially life-threatening, whereas misidentifying a *Normal* conclusion as *Non-normal* will only incur extra cost to conduct further tests. To investigate this point, we have compared the probabilities $Pr(C = NNor|A_Q = Yes)$ with $Pr(C = NNor|A_Q = No)$ and used the Z-test [22] to check whether the difference between the two resulting Bernoulli distributions is statistically significant. The procedure is as follows: Given a question Q , let $T_{A_Q=Yes}$ be the total number of participants answering *Yes* while $T_{A_Q=No}$ denotes the total number of participants answering *No* to the question. The Z-statistic for the probabilities $Pr(C = NNor|A_Q = Yes)$ and $Pr(C = NNor|A_Q = No)$ is calculated as:

$$Z = \frac{Pr(C=NNor|A_Q=Yes) - Pr(C=NNor|A_Q=No)}{\sqrt{p(1-p)(1/T_{A_Q=Yes} + 1/T_{A_Q=No})}},$$

Table 4. Classification results from the WEKA implementation of naïve Bayes, random forests (RF) and support vector machine (SVM), on the original (biased) dataset

Classifier	Accuracy for Normal class	Accuracy for Non-normal class	Overall Accuracy	Precision	Recall	F-measure
Naïve Bayes	0.968	0.098	0.742	0.522	0.098	0.166
RF	0.905	0.115	0.70	0.298	0.115	0.166
SVM	0.968	0.098	0.742	0.522	0.098	0.166

$$\text{where, } p = \frac{T_{A_Q=No} * Pr(C=NNor | A_Q=Yes) + T_{A_Q=Yes} * Pr(C=NNor | A_Q=No)}{T_{A_Q=Yes} + T_{A_Q=No}}.$$

For a two-sided test, if the value of the Z-statistic is greater than 1.96 or smaller than -1.96, the difference between the two probabilities is considered statistically significant with 95% confidence (p-value ≤ 0.05).

We also examined the (lack-of) association between affirmative answers to *combinations of questions* and the *Non-normal* conclusion. The details are not described here due to space limitation and will be included in an extended version of this paper.

4 Results

As mentioned in Section 1, as a baseline, we set out to classify the dataset using traditional machine learning methods: naïve Bayes, random forests, and support vector machine. The goal was to assign the athletes into the correct adjudicated class (i.e., predict the ECG/Echo conclusion), based on their respective answers to the questions shown in Table 2. However, *all three classifiers performed poorly* for the *Non-normal* class, as evaluated using 10-fold cross validation. The classification Accuracy, Precision, Recall and F-measure for the three methods when applied to the original (biased) dataset are shown in Table 4. For the *Normal* class, the naïve Bayes, the random forest and the SVM classifiers correctly classified 96.8%, 90.5% and 96.8% instances, respectively, but their performance for the *Non-normal* class is extremely poor. As noted before, the performance over the *Non-normal* class is very important because misclassifying an athlete with an abnormal heart condition as *Normal* is unacceptable in a pre-screening process. We note that the poor performance of the classification for *Non-normal* class may be attributed to the bias in the dataset, which can lead the classifier to assign most of the instances to the majority class. To correct for this, we have used sub-sampling for balancing the set; Table 5 shows the classification results for the balanced datasets, averaged over 5 random sub-samples. Correcting for the imbalance in the dataset indeed improved significantly the classification results for instances of the *Non-normal* class (the Recall in particular), but still, about 50% of the *Non-normal* cases are misclassified as *Normal* by naïve Bayes and 36% are misclassified as *Normal* by random forests. Similarly the SVM classifier misclassifies 45% of the *Non-normal* cases as *Normal*. Moreover, the vast majority of the *Normal* cases (more than 50%, for all three classifiers) have been classified as *Non-normal*. Notably, such a low level of performance is close to the classification level expected at random.

As discussed in Section 3.2, to pursue the information-content based analysis of each question, we calculated the information gain per question. The information gain associated with questions *Q1-Q12* ranges between 0.001-0.003 and for *Q13* it is 0.008. Clearly, the information gain for all of the questions is very low, the highest b-

Table 5. Classification results from the WEKA implementation of naïve Bayes, random forests (RF) and support vector machine (SVM) on the balanced dataset

Classifier	Accuracy for Normal class	Accuracy for Non-normal class	Overall Accuracy	Precision	Recall	F-measure
Naïve Bayes	0.443	0.508	0.475	0.477	0.508	0.492
RF	0.467	0.639	0.553	0.545	0.639	0.589
SVM	0.459	0.549	0.504	0.504	0.549	0.525

eing only 0.008 for question *Q13*, which is the result of the AHA-recommended physical exam. As a point of comparison, in a hypothetical case in which even just 70% of the *Yes* answers to question *Q13* would corresponded to a *Non-normal* conclusion, the information gain would have been 0.106, which is significantly higher than any of the gains associated with the questions. This very low information content of each question explains the poor classification results, especially the close-to-random classification performance over the balanced dataset.

To further analyze whether positive answers to the questions have higher probability of corresponding to *Non-normal* conclusion than negative answers, we have compared the probabilities $Pr(C = NNor|A_Q = Yes)$ and $Pr(C = NNor|A_Q = No)$. The histogram in Figure 1 shows for each question the conditional probability of the conclusion being *Non-normal* given that the answer to the question is *Yes*, side-by-side with the conditional probability of a *Non-normal* conclusion, when the answer to the same question is *No*. We observe that for seven of the questions (*Q3*, *Q4*, *Q5*, *Q9*, *Q11*, *Q12* and *Q13*), the conditional probability $Pr(C = NNor|A_Q = Yes)$ is indeed somewhat higher than the conditional probability $Pr(C = NNor|A_Q = No)$. However, for six of the questions, *Q1*, *Q2*, *Q6*, *Q7*, *Q8*, and *Q10*, the probability of a *Non-normal* adjudication is actually *higher* when the answer is negative than when the answer is positive. We used the Z-test to verify whether these differences are statistically significant, and found that only for *Q13* (the physical exam), the difference is statistically significant with a p-value of 0.016. Thus the only item in the questionnaire for which a positive answer is marginally predictive of a *Non-normal* conclusion, is the physical examination (*Q13*). However, even in this case, the number of false negatives (i.e. the number of *Non-normals* that are left undetected) is 94 out of a total of 110 *Non-normals*, which is very high.

All of the results described above demonstrate that relying on normal findings from the physical examination (*Q13*), and on negative answers to questions *Q1-Q12* in the AHA questionnaire as a way to assess whether athletes can safely participate in competitive activities *leads to a high rate of false negatives*. That is, athletes with potential heart abnormalities (identified by a cardiologist through ECG and Echo tests) are very likely to be pre-screened as *Normal*, and not be referred for further examination. This is clearly an undesirable scenario in a pre-screening process.

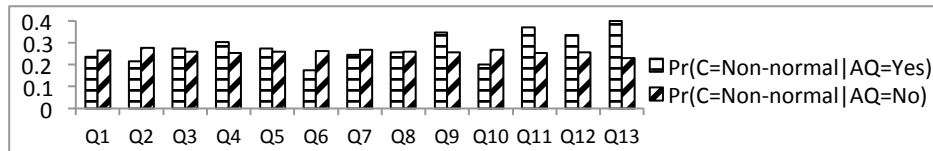


Fig. 1. Conditional probability of adjudications being *Non-normal* when the answer to each question is *Yes* vs. *No*

5 Conclusion

We set out to build a classifier that could predict potential abnormalities in young athletes' heart-condition, using data from close to 500 athletes who were examined using the AHA-based 12-element screening procedure. The ground truth used for potential abnormality was determined by an expert cardiologist based on Electro- and Echo-cardiogram tests, which are not included in the AHA screening procedure.

The poor performance of several well-studied machine-learning classifiers, (and particularly the close-to-random classification performance measured on the balanced dataset), when using all the elements in the questionnaire as attributes, lead us to conduct an in-depth study of the data and the questions. We aimed to determine each element's ability (or there lack-of) to identify abnormality. Underlying this part of the study was the expectation that the classifiers' performance may be improved by using the most informative subset of questions as attributes. However, surprisingly, our results show that in terms of information content, none of the elements included in the questionnaire contributes significant information about the findings obtained through traditional ECG and Echo-based tests. *As such, improvement in the classification results is not attainable using any subset of the questions as attributes. Through the use of machine-learning and statistical methods, we identified that the culprit is in the screening procedure itself.* Further analysis of the respective conditional probabilities through statistical tests, indicates that an abnormal physical examination (*Q13*) is the only item within the questionnaire that is even associated with a statistically-significantly higher probability of a *Non-normal* ECG/Echo than a normal physical examination. But even this item still gives rise to many false negatives.

Thus, the results of this study are highly significant, as they strongly suggest that the 12-element procedure advocated by the American Heart Association for pre-participation screening of young athletes is not correlated with or predictive of the outcome obtained by a cardiologist using standard ECG and Echo tests.

Pragmatically speaking, the conclusion from this study implies that ECG (and possibly Echo) should be considered for screening athletes in the Unites States. Future research following the machine-learning and informatics-driven approach as used in this study will examine whether using one or more of the cardiovascular tests such as electrocardiogram or echocardiogram together with any combination of all or some of the AHA-based questions may improve the efficacy of pre-participation screening.

Acknowledgments. This work was partially supported by HS's NSERC Discovery Award #298292-2009, NSERC DAS #380478-2009, CFI New Opportunities Award 10437, and Ontario's Early Researcher Award #ER07-04-085, and by TA's grant HL 098046 from the National Institutes of Health.

References

1. Corrado, D., et al.: Cardiovascular Pre-Participation Screening of Young Competitive Athletes for Prevention of Sudden Death: Proposal for A Common European Protocol. Consensus statement of the study grp. of Sport Cardiology, of the wrk. grp. of Cardiac Rehabilitation and Exercise Physiology and the wrk. grp. of Myocardial and Pericardial Disease of the European Society of Cardiology. *European Heart J.* 26(5), 516–24 (2005)

2. Maron, B.J.: Sudden Death in Young Athletes. *New England J. of Medicine* 349(11), 1064–75 (2003)
3. Pigozzi, F., Rizzo, M.: Sudden Death in Competitive Athletes. *Clinics in Sports Medicine* 27(1), 153–81 (2008)
4. Wever-Pinzon, O.E., et al.: Sudden Cardiac Death in Young Competitive Athletes Due to Genetic Cardiac Abnormalities. *Anadolu Kardiyol Derg* 9(Suppl 2), 17–23 (2009)
5. Corrado, D., et al.: Screening for Hypertrophic Cardiomyopathy in Young Athletes. *New England J. of Medicine* 339(6), 364–69 (1998)
6. Maron, B.J., et al.: Cardiovascular Preparticipation Screening of Competitive Athletes: A Statement for Health Professionals From the Sudden Death Committee (Clinical Cardiology) and Congenital Cardiac Defects Committee (Cardiovascular Disease in the Young). *Circulation* 94(4), 850–56 (1996)
7. Glover, D.W., Maron, B.J.: Evolution in the Process of Screening United States High School Student-athletes for Cardiovascular Disease. *American J. of Cardiology* 100(11), 1709–12 (2007)
8. Maron, B.J., et al.: Recommendations and Considerations Related to Preparticipation Screening for Cardiovascular Abnormalities in Competitive Athletes: 2007 Update: A Scientific Statement from the American Heart Association Council on Nutrition, Physical Activity, and Metabol. *Circulation* 115(12), 1643–55 (2007)
9. Kanagalingam, J., et al.: Efficacy of the American Heart Association Questionnaire in Identifying Electrocardiographic and Echocardiographic Abnormalities in Young Athletes During Community-based Screening. *Circulation* 122(21), A19765 (2010)
10. Mitchell, T.M.: *Machine Learning*. McGraw-Hill (1997)
11. Melgani, F., Bazi, Y.: Classification of Electrocardiogram Signals with Support Vector Machines and Particle Swarm Optimization. *IEEE Trans. on Information Technology in Biomedicine* 12(5), 667–77 (2008)
12. Osowski, S., Hoai, L.T., Markiewicz, T.: Support Vector Machine-Based Expert System for Reliable Heartbeat Recognition. *IEEE Trans. on Biomedical Eng.* 51(4), 582–89 (2004)
13. Chazal, P. D., et al.: Automatic Classification of Heartbeats using ECG Morphology and Heartbeat Interval Features. *IEEE Trans. on Biomedical Eng.* 51(7), 1196–1206 (2004)
14. Yu, S., Chou, K.: Integration of Independent Component Analysis and Neural Networks for ECG Beat Classification. *Expert Systems with Applications* 34(4), 2841–2846 (2008)
15. Qu, L., et al.: A Naïve Bayes Classifier for Differential Diagnosis of Long QT Syndrome in Children. In: *Int. Conf. on Bioinformatics and Biomedicine*, pp. 433–37 (2010)
16. Akay, M.F.: Support Vector Machines Combined with Feature Selection for Breast Cancer Diagnosis. *Expert Systems with Applications* 36(2), 3240–47 (2009)
17. Chhatwal, J., et al.: A Logistic Regression Model Based on the National Mammography Database Format to Aid Breast Cancer Diagnosis. *American J. of Roentgenology* 192(4), 1117–27 (2009)
18. Statnikov, A., Wang, L.: A Comprehensive Comparison of Random Forests and Support Vector Machines for Microarray-Based Cancer Classification. *BMC Bioinformatics* 9(1), 319 (2008)
19. Breiman, L.: Random forests. *Machine learning* 45(1), 5-32 (2001)
20. Cortes, C., Vapnik, V.: Support-vector Networks. *Machine learning* 20(3), 273–297 (1995)
21. Hall, M., et al.: The WEKA Data Mining Software: an Update. *ACM SIGKDD Explorations Newsletter* 11(1), 10-18 (2009)
22. Walpole, R., et al.: *Probability and Statistics for Engineers and Scientists*. Prentice Hall (2002)